

---

## GELATO@UNSW VM WORK

Slide 1



Ian Wienand, Paul Davies, Lucy Chubb

This work supported by UNSW and HP through the Gelato Federation

---

## ABOUT THE TALK

- Work in progress talk
- Please talk to us!

Slide 2

### Topics:

- Superpages
  - Long Format VHPT
  - Page Table Abstraction
- 

---

## SUPERPAGES

### What?:

- Use page size > base page size
- HugeTLB
  - ✗ Programmer Visible
  - ✗ Confusing
  - ✓ Suits legacy architectures (x86/64, Power)

Slide 3

### Why?:

- More coverage from TLB
  - Transparency is a worthwhile, unsolved problem
  - Itanium provides wide range of interesting MMU options
    - **With LVHPT!**
- 

---

## OUR APPROACH

### Cluster individual pages:

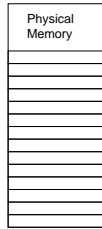
Slide 4

- ① Tag entries with size when created - `mmap(MAP_ANONYMOUS)`
  - ② On fault, check page size flag
  - ③ Backtrack and **allocate physical superpage**
  - ④ Point each sub-page of superpage to actual physical page
-

Slide 5

## SUPERPAGE OVERVIEW

### PTE Entries



P	Physical	Size
	0x00000000	0
	0x00000000	0
	0x00000000	8
	0x00000000	8
	0x00000000	8
	0x00000000	8
	0x00000000	8
	0x00000000	8
	0x00000000	8
	0x00000000	0
	0x00000000	0

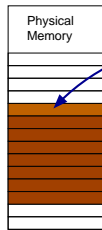
p = mmap(512K, MAP\_ANON)



Slide 6

## SUPERPAGE OVERVIEW

### PTE Entries



P	Physical	Size
	0x00000000	0
	0x00000000	0
	0x10010000	8
	0x10020000	8
	0x10030000	8
	0x10040000	8
	0x10050000	8
	0x10060000	8
	0x10070000	8
	0x10080000	8
	0x00000000	0
	0x00000000	0

p = mmap(512K, MAP\_ANON)

p[1] = 100;

Page Fault

Allocate Memory

Retry Access

To TLB

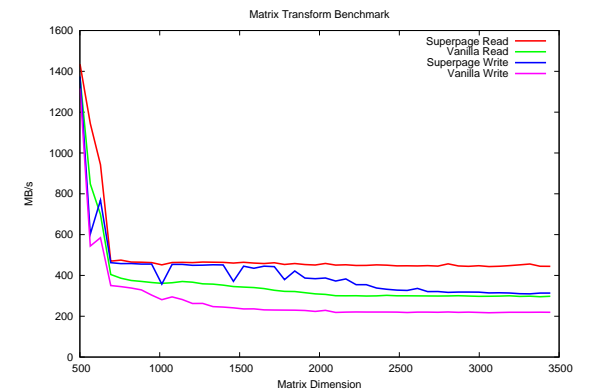
No more page faults!

## ISSUES

- Work with Itanium features
- Page Tables
  - Double Size PTE
  - Mapping across PMD's
  - Better approaches
- Fragmentation
  - External Fragmentation (Ozlabs, Rice)
  - Internal Fragmentation (Rice)
- Promotion/Demotion
  - Competitive Algorithms

Slide 7

## MATRIX TRANSFORMATION



Slide 8

## DODGY RESULTS

Slide 9

- lmbench - no massive differences
- SPECint - mcf shows good improvement
- Lee's benchmark?

## VIRTUALLY HASHED PAGE TABLE

Problem:

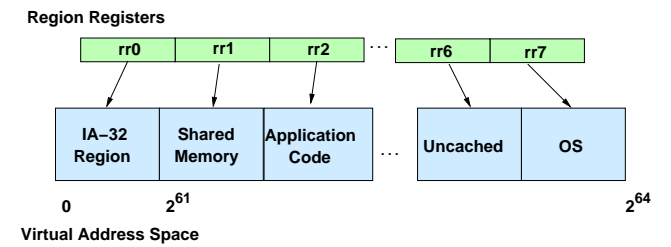
- TLB is **fast** but **small**
- Walking page tables via software is **slow**

Solution:

- Keep a "cache" of translations in RAM
- Let the hardware do the walking!
- The **VHPT**

## FAST VM OVERVIEW - REGIONS

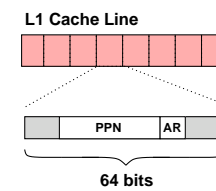
Slide 11



## TWO TYPES OF VHPT

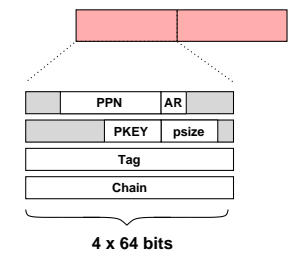
Slide 12

### Short Format



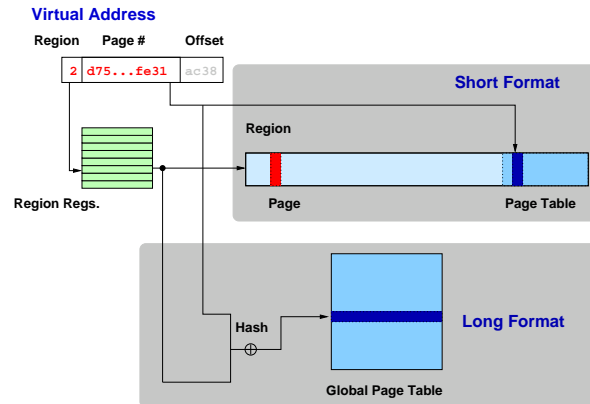
8 entries per cache line  
Page Size from Region

### Long Format



2 entries per L1 cache line  
Page tables NOT 4 times bigger

## VHPT DIFFERENCES



Slide 13

## TWO TYPES OF VHPT

- **Short** : Cache Friendly, TLB Unfriendly
  - More entries in a cache line
  - Extra TLB entry per page of PTEs
- **Long** : Cache Unfriendly, TLB Friendly
  - Less entries per cache line
  - Global page table pinned with single entry

Slide 14

Why long?:

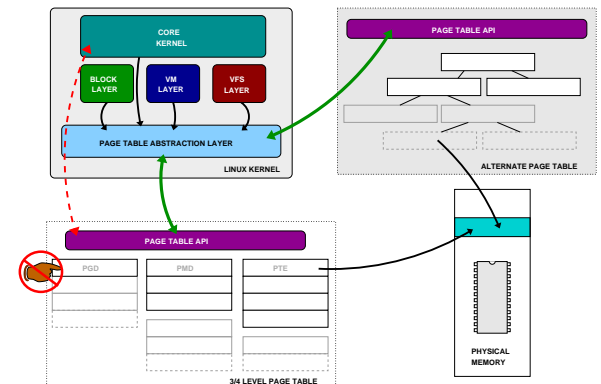
- ✓ Entries can hold page size, protection keys
- ✗ Entries are twice as big.
- Would like to talk to hackers and users.

## PAGE TABLE ABSTRACTION

Slide 15

- Problems with current implementations:
- 3 (4? 5? X?) level page table sub-optimal
  - Main architectures have differing needs

## SOLUTION? PAGE TABLE ABSTRACTION



Slide 16

---

Issues:

Slide 17

- ✓ Performance - not bad!
- ✓ Hardware walked page tables - we have ideas
- ✓ Alternate Implementations - GPT
- ✗ Convincing Linux community to rip apart VM - hard

---

WHERE DO YOU GET IT?

<http://www.gelato.unsw.edu.au/IA64Wiki>

or

Google for "ia64 wiki"

QUESTIONS?

---

Slide 18

QUESTIONS?